
Predicting Shot-Creating Passes in Soccer Using Transformer Models: A Sequence-Based Analysis of Leicester City’s 2015–16 Season

Kejian Zhu
University of Michigan
zhukj@umich.edu

Abstract

1 Understanding how danger emerges in soccer requires modeling the sequential
2 structure of match events rather than treating actions in isolation. This paper investi-
3 gates whether a Transformer encoder can identify pass sequences that precede
4 dangerous outcomes, defined as shots taken within ten events following a pass with
5 expected goals (xG) of at least 0.05. Using complete event data from Leicester
6 City’s 2015–16 Premier League campaign, the model processes six-pass windows
7 and predicts whether the final pass will lead to a shot. A logistic regression classi-
8 fier serves as a baseline and is evaluated on the same sequence-level test set. The
9 logistic model achieves an AUC of 0.887 and AP of 0.308, while the Transformer
10 achieves an AUC of 0.887 and a substantially higher AP of 0.410. Although
11 the global discrimination (AUC) is similar, the Transformer exhibits significantly
12 improved ranking of rare dangerous events, suggesting that sequence context and
13 player-interaction embeddings meaningfully enhance predictive precision. Inter-
14 pretability analyses demonstrate that the model captures tactical patterns consistent
15 with Leicester City’s transitional attacking style. These findings highlight the value
16 of sequence-based deep learning approaches for soccer event prediction.

17 **1 Introduction**

18 Soccer is a sequential, context-dependent game in which the value of each action depends strongly
19 on the preceding flow of play. A single pass can be safe or extremely dangerous depending on prior
20 movements, field position, and tactical structure. Traditional soccer analytics often models actions
21 in isolation, focusing on single events such as shots, passes, or duels, and assigning them values
22 through possession-value frameworks like Expected Threat (xT) (1). While these approaches have
23 been successful in quantifying the contribution of individual actions, they typically assume that threat
24 is primarily a function of location and immediate context rather than longer attacking sequences.

25 The 2015–16 Leicester City team provides a compelling case study for sequence-based modeling.
26 Leicester’s title-winning campaign was built on rapid transitions, vertical attacks, and coordinated
27 interplay among players such as Jamie Vardy, Riyad Mahrez, N’Golo Kanté, Marc Albrighton, and
28 Christian Fuchs. Many of their most dangerous situations emerged not from isolated high-value
29 passes but from multi-pass sequences that stretched the opposition and created space before the final
30 action. Capturing these dynamics requires models that can reason about temporal context and player
31 interaction patterns across multiple events.

32 This paper investigates whether a Transformer-based sequence model can predict dangerous passes
33 more effectively than a simple logistic regression baseline. Using event-level data from Leicester
34 City’s 2015–16 Premier League matches, we construct short sequences of six consecutive Leicester
35 passes and predict whether the final pass in the window is followed by a Leicester shot within the
36 next ten events, with expected goals (xG) at least 0.05. The Transformer processes the entire six-pass

37 sequence, while the logistic regression baseline uses only the numeric features of the final pass. Both
38 models are trained and evaluated on the same sequence-level train, validation, and test splits, allowing
39 for a fair comparison.

40 The main contributions of this work are threefold. First, we introduce a sequence-based prediction
41 task focused on identifying shot-creating passes in a real-world season for a specific team. Second, we
42 implement and evaluate a Transformer architecture that combines spatial features, temporal ordering,
43 and player identity embeddings, and we compare it against a logistic regression baseline trained on
44 the same sequence endpoints. Third, we provide an interpretability analysis of the learned model,
45 illustrating how it captures tactical patterns consistent with Leicester City’s transitional attacking
46 style.

47 2 Related Work

48 Soccer analytics has increasingly turned to event and tracking data to quantify performance, evaluate
49 player contributions, and model match dynamics. Possession-value models such as Expected Threat
50 (xT) assign values to locations on the pitch based on the probability that a possession will eventually
51 result in a goal (1; 2). Recent extensions like Dynamic Expected Threat (DxT) incorporate player
52 positions to make action evaluation more context-sensitive (3). These frameworks typically operate
53 at the level of single actions, updating threat values as the ball moves, but do not explicitly model
54 sequences of actions as inputs to a predictive model.

55 Several studies have explored machine learning approaches for predicting outcomes of possessions or
56 phases of play using event data. Stival et al. use a machine learning pipeline to predict whether a
57 team successfully enters the attacking third within a few seconds of gaining possession (4). Other
58 work has used event data from providers such as StatsBomb and Wyscout to model match outcomes
59 or to quantify player impact across matches (5; 6). These approaches focus on higher-level outcomes
60 (e.g., match result, attack success) rather than directly modeling the probability that a specific pass
61 will lead to a shot.

62 Deep learning methods and sequence models have recently been applied to soccer event streams.
63 Simpson et al. introduced Seq2Event, a framework that uses recurrent or Transformer encoders to
64 predict the next event and its attributes from sequences of match events (7). Transformer-based
65 spatio-temporal point process models have been proposed to jointly predict the timing, location, and
66 type of future events (8). Other recent work has used Transformers to assess player performance
67 and to provide a versatile framework for analyzing soccer actions (9). Beyond soccer, sequence
68 models, including Transformers, have been successfully applied to event prediction in other sports
69 and domains (10; 11).

70 Passing sequences themselves have been studied using network and pattern-mining approaches.
71 McCarthy et al. analyze passing sequences to infer the style of play and relate them to goal-scoring
72 opportunities (12). Clustering and motif-based methods have been used to identify line-breaking
73 passes and to uncover tactical structures that create or reduce threat (13). These approaches often
74 rely on handcrafted features or unsupervised clustering, whereas our work trains a supervised model
75 end-to-end to predict whether a sequence will culminate in a dangerous shot.

76 Our study is closest in spirit to Seq2Event and Transformer-based soccer models, but differs in its
77 focus on a specific, interpretable outcome: whether a six-pass sequence leads to a shot within the next
78 ten events. By restricting the analysis to a single season for a single team, we can connect the learned
79 patterns to known tactical structures rather than aiming for a fully general model across leagues.

80 3 Method

81 3.1 Data and Preprocessing

82 We use event-level data from StatsBomb Open Data for the 2015–16 Premier League season (14).
83 All matches involving Leicester City are retrieved via the StatsBombPy API. For each match, we
84 access the full sequence of events, including passes, shots, duels, and other actions, and attach the
85 match identifier.

86 The preprocessing pipeline focuses on Leicester City passes. We first filter the event log to retain
87 only Leicester events and then further restrict to events labeled as passes. For each pass, we extract
88 the starting and ending locations on the pitch and compute engineered spatial features, including
89 normalized coordinates on a 120-by-80 pitch, pass length, and pass angle. We convert boolean flags
90 such as under pressure and counterpress into binary indicators and encode categorical variables,
91 including pass height, body part, pass type, and play pattern, using numeric category codes. Player
92 and pass recipient identities are encoded as integer indices, with a special index reserved for missing
93 recipients to ensure that all passes have a valid receiver code.

94 To associate each pass with a danger label, we sort all events in each match by period, minute, second,
95 and internal index, and assign an event index. A pass is labeled as dangerous if, within the next ten
96 events in the same match, there exists a Leicester shot with expected goals (xG) at least 0.05. This
97 threshold filters out very low probability attempts while retaining shots that represent meaningful
98 attacking outcomes. The label is binary and does not distinguish between different shot qualities
99 beyond the threshold.

100 Finally, we drop passes with missing key features or invalid player identifiers and standardize numeric
101 features using a global StandardScaler. This cleaned dataframe provides the basis for both the logistic
102 regression baseline and the sequence construction.

103 3.2 Sequence Construction

104 The task is to predict whether the final pass in a short attacking sequence will lead to a shot in the
105 near future. To create sequences, we group passes by match and sort them chronologically within
106 each match. For each match, we slide a fixed-length window of six consecutive Leicester passes. For
107 a given window, the model input consists of the numeric features of each pass, the passer and receiver
108 indices, and the positional index (0 to 5) in the sequence. The target label is the danger label of the
109 final pass in the window.

110 Only matches with at least six Leicester passes contribute sequences. For each valid window, we
111 obtain a tensor of shape $[L, F]$ for numeric features, where $L = 6$ and F is the number of numeric
112 features, as well as vectors of passer and receiver indices of length L . The corresponding label is a
113 scalar indicating whether the final pass leads to a dangerous shot.

114 We construct the full set of sequences and then split them into training, validation, and test subsets
115 using stratified sampling on the labels. This ensures that each split maintains a similar proportion
116 of dangerous and non-dangerous sequences, which is important because dangerous sequences are
117 relatively rare.

118 3.3 Models

119 **Logistic Regression Baseline.** The baseline model is a logistic regression classifier trained on
120 the same sequence dataset but using only the numeric features of the final pass in each six-pass
121 window. For a given sequence, we extract the feature vector of the last pass and use it as input to the
122 logistic regression. This ensures that the logistic model operates on the same samples and labels as
123 the Transformer, making the comparison fair. The classifier is trained with class-balanced weights
124 and a maximum of 1000 iterations.

125 **Transformer Sequence Model.** The main model is a Transformer encoder that processes the full
126 six-pass sequence. Numeric features for each pass are projected into a latent space via a linear layer.
127 Passer and receiver indices are mapped to learned embeddings, which are also projected to the same
128 latent dimension. Positional embeddings indicate the temporal order of passes within the window.
129 The sum of the projected numeric features, passer embeddings, receiver embeddings, and positional
130 embeddings yields the input to the Transformer encoder.

131 The encoder comprises several layers of multi-head self-attention and feedforward sublayers, with
132 layer normalization and dropout for regularization. The output is a contextualized representation
133 for each pass in the sequence. We extract the hidden state corresponding to the final pass and feed
134 it into a small feedforward network that outputs a scalar logit. The model is trained using binary
135 cross-entropy with logits, with an additional positive class weight to address label imbalance. The

136 Adam optimizer is used with a fixed learning rate, and the best model is selected based on validation
137 AUC.

138 **4 Experiment**

139 **4.1 Experimental Setup**

140 After constructing all valid sequences, we perform a stratified split into training, validation, and test
141 subsets. The logistic regression baseline and the Transformer are both trained using the training set,
142 with hyperparameters tuned informally using the validation set. Final performance is reported on the
143 held-out test set only.

144 We evaluate both models using two metrics: the area under the ROC curve (AUC) and average
145 precision (AP). AUC measures the overall ability of the model to rank dangerous sequences above
146 non-dangerous ones across all thresholds. AP summarizes the precision–recall trade-off and is
147 particularly informative in imbalanced settings where positive examples are rare. In the context of
148 football analytics, AP is arguably more aligned with practical use cases, since analysts are often
149 interested in the reliability of the top-ranked predicted dangerous passes.

150 **4.2 Results**

151 On the sequence-level test set, the logistic regression baseline achieves an AUC of 0.887 and an AP
152 of 0.308. These results show that the numeric features of the final pass alone contain substantial
153 information about whether a shot will occur in the near future. In particular, spatial location, pass
154 length, angle, and contextual attributes such as pass type and play pattern already allow a linear model
155 to distinguish many dangerous from non-dangerous actions.

156 The Transformer model, trained on the full six-pass sequence with player and receiver embeddings,
157 obtains an AUC of approximately 0.88 and an AP of 0.41 on the same test set. The AUC is very
158 similar to that of the logistic regression baseline, suggesting that both models achieve comparable
159 global discrimination when ranking all sequences. However, the AP improves noticeably from 0.308
160 to around 0.41. This indicates that the Transformer is better at concentrating true dangerous sequences
161 near the top of its ranking, which is particularly valuable when analysts focus on the most dangerous
162 predicted events.

163 The fact that AUC remains similar while AP improves is consistent with the idea that sequences add
164 subtle but important information that is most useful in the high-risk region. The logistic model is
165 competitive when evaluating all thresholds, but the Transformer provides a more refined ranking of
166 the most threatening passes.

167 **4.3 Interpretability and Tactical Insights**

168 Although the primary goal is predictive performance, the sequence model also yields insights into
169 Leicester City’s attacking structure. By examining the learned representations and scores, we observe
170 that dangerous sequences often involve repeated interactions among well-known attacking players.
171 Passer–receiver pairs such as Albrighton to Huth or Morgan, Fuchs to central defenders or Mahrez,
172 and combinations involving Kanté, Mahrez, Okazaki, and Vardy appear frequently in high-probability
173 sequences. These patterns align with established descriptions of Leicester’s transitional style, in
174 which the team built attacks via wide distribution from full-backs and deep players before accelerating
175 through Mahrez and Vardy.

176 Sequence-level importance scores, computed by comparing intermediate Transformer hidden states
177 to the representation of the final pass, suggest that early passes in a six-pass window often receive
178 substantial weight. This supports the notion that danger is not solely determined by the final action,
179 but by how earlier passes shape the defensive structure and field position. Spatial plots of final pass
180 locations, colored by predicted probability, show dense clusters in wide and half-space regions near
181 the penalty area, consistent with cutback and crossing zones that are known to be high-value in
182 practice.

183 **5 Conclusion and Discussion**

184 This paper presented a sequence-based approach to predicting shot-creating passes in soccer using
185 a Transformer model applied to Leicester City’s 2015–16 Premier League season. By constructing
186 six-pass windows and predicting whether the final pass leads to a shot within the next ten events, we
187 framed the problem as a supervised sequence classification task. A logistic regression baseline trained
188 on the same sequence endpoints achieved strong AUC but relatively modest average precision. The
189 Transformer matched the baseline in AUC while significantly improving AP, indicating that sequential
190 context and player-interaction embeddings help the model prioritize truly dangerous sequences near
191 the top of the ranking.

192 From a methodological standpoint, the results show that even in a relatively small, team-specific
193 dataset, deep sequence models can provide meaningful gains over simple baselines when evaluated
194 with ranking-focused metrics. At the same time, the strong performance of the logistic regression
195 model cautions against overstating the advantage of deep models; much of the signal remains encoded
196 in spatial and contextual features of the final pass itself. The primary value of the Transformer appears
197 in its ability to exploit temporal dependencies and repeated interaction patterns that are not fully
198 captured by single-pass features.

199 For soccer practitioners, this framework suggests a way to identify and quantify multi-pass sequences
200 that are likely to lead to shots, with clear tactical interpretations. Analysts could apply similar models
201 to other teams or seasons to compare attacking structures, evaluate how player departures or signings
202 change dangerous combinations, or simulate how different sequence patterns affect the likelihood of
203 shot creation.

204 Future work could extend this study in several directions. Incorporating tracking data would allow
205 the model to account for defender positions and off-ball movements, likely improving both predictive
206 accuracy and tactical interpretability. Training on multi-team or multi-league datasets would test
207 the generality of the learned patterns and potentially support transfer learning applications. Finally,
208 integrating counterfactual or reinforcement learning frameworks could enable simulations in which
209 alternative passing choices within a sequence are evaluated in terms of their effect on the probability
210 of generating a dangerous shot.

211 Overall, this case study illustrates how modern sequence models, and Transformers in particular,
212 can be adapted to soccer analytics to move beyond single-action evaluation and toward a richer
213 understanding of how coordinated sequences of events create threat.

214 The code is available at: <https://github.com/zhukejianbot/LeicesterPass>

215 **References**

- 216 [1] Karun Singh. Introducing Expected Threat (xT): Modelling team behaviour in possession to
217 gain a deeper understanding of buildup play. Blog post, 2018. Available at <https://karun.in/blog/expected-threat.html>.
- 219 [2] Footballitics. Expected Threat (xT): the best offensive metric so far. Blog post, 2022.
- 220 [3] Karim Hassani et al. Dynamic Expected Threat (DxT) Model: Addressing the Deficit of Realism
221 in Football Action Evaluation. *Applied Sciences*, 15(8), 2025.
- 222 [4] Leandro Stival et al. Using machine learning pipeline to predict entry into the attack zone in
223 soccer. *PLOS ONE*, 18(2): e0265372, 2023.
- 224 [5] Peter Hassard et al. Predicting football match outcomes using event data and machine learning
225 algorithms. Technical report, 2024.
- 226 [6] Tahmeed Tureen and Sigrid Olthof. Estimated Player Impact (EPI): Quantifying the Effects of
227 Individual Players on Football Actions. StatsBomb Conference Paper, 2022.
- 228 [7] Ian Simpson et al. Seq2Event: Learning the Language of Soccer Using Transformer-based
229 Event Representations. In *Proceedings of KDD*, 2022.
- 230 [8] Jing Li et al. Transformer-based Neural Marked Spatio-Temporal Point Process for Football
231 Events. *Applied Intelligence*, 2025.

- 232 [9] Andrey Rovshitz et al. SoccerTransformer: A Transformer-based Framework for Versatile
233 Analysis of Soccer Players. In *Proceedings of MLSA*, 2024.
- 234 [10] Jian Gao et al. Predicting Sport Event Outcomes Using Deep Learning. *Journal of Sports*
235 *Analytics*, 2025.
- 236 [11] Kevin McGuigan et al. A hard-hitting evaluation of deep learning approaches to predicting
237 tackle events in the NFL. *Vision Science Letters*, 2024.
- 238 [12] Conor McCarthy et al. Analyzing Passing Sequences for the Prediction of Goal-Scoring
239 Opportunities. In *Workshop on Machine Learning and Data Mining for Sports Analytics*
240 (*MLSA*), 2022.
- 241 [13] Marc Lamberts et al. Uncovering Tactical Line-Breaking Passes with Clustering. arXiv preprint,
242 2025.
- 243 [14] StatsBomb. StatsBomb Open Data. GitHub repository, available at <https://github.com/statsbomb/open-data>.